

# Processamento de Cadeias de Caracteres

Prof. José J. Camata

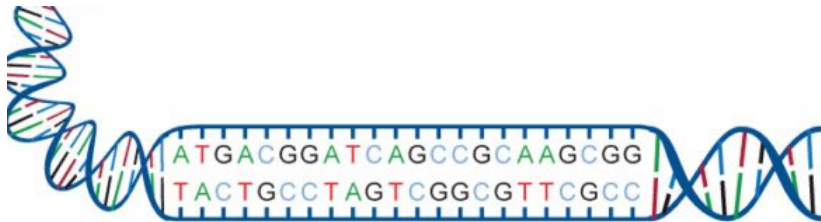
Prof.a. Bárbara Quintela

[camata@ice.ufjf.br](mailto:camata@ice.ufjf.br)

[barbara@ice.ufjf.br](mailto:barbara@ice.ufjf.br)

# Cadeias de Caracteres

- Uma cadeia é uma sequência linear de caracteres, tipicamente podendo ser muito longa
- Cadeias são centrais em sistemas de processamento de textos, filtro de spam, recuperação de informação, detecção de plágio, estudo de sequências de DNA em biologia computacional, etc.
- Essas cadeias podem ser bem grandes e algoritmos eficientes são necessários para manipulá-los.



# Cadeias de Caracteres

- Abordaremos algumas questões relacionadas com cadeias de caracteres nos seguintes tópicos
  - **Casamento de padrões em cadeia de caracteres (String matching)**
  - **Codificação e Compressão**

# Cadeias de Caracteres

## ➤ Notação:

- $T$ , usado para representar um texto, ou seja, uma sequência de caracteres ou letras.
- $P$  é um sequência de caracteres que se deseja procurar em  $T$
- $|T|$  significa o comprimento de  $T$
- $T_j$  é o caracter na posição  $j$ .
- $T[i, \dots, j]$  é um sequência de  $T$  que começa na posição  $i$  e termina em  $j$ .
- Os primeiros caracteres do padrão  $P$  e o texto  $T$  estão na posição 1.

## ➤ Exemplo:

- $T = \text{"ABACAABACC**ABAC**ABAABB"}$
- $|T| = 20$
- $P = \text{"ABACAB"}$  corresponde a trecho  $T[11, \dots, 16]$

# Casamento de Padrão

- Uma operação fundamental sobre cadeias é o casamento de padrão:
  - consiste encontrar uma cópia exata do padrão  $P$  de tamanho  $m = |P|$  no texto  $T$  de tamanho  $n = |T|$ .
  - O problema de casamento de padrão pode ser visto também como um problema de busca com o padrão sendo a chave.
- Duas soluções serão estudadas:
  - **Força Bruta:** simples mas pouco eficiente. Tem complexidade  $O(nm)$
  - **Algoritmo de Knuth, Morris e Pratt:** mais elaborado e com complexidade  $O(n)$ .

## Nota histórica

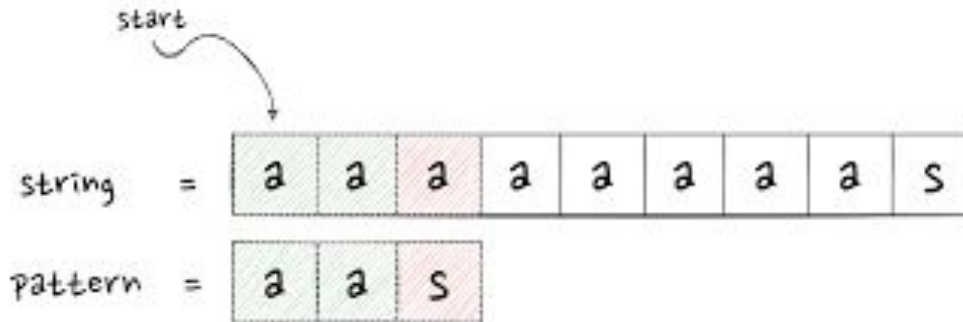
- Em 1970, S.A. Cook provou um resultado teórico sobre um tipo particular de autômato que implicava na existência de um algoritmo de casamento de padrão com tempo proporcional a  $M + N$  no pior caso.
- D. E. Knuth e V. R. Pratt seguindo a construção que Cook usara na demonstração do seu teorema obtiveram um algoritmo relativamente simples e prático.
- Ocorreu também que J. H. Morris descobriu praticamente o mesmo algoritmo como solução de um problema de edição de texto
- Os três cientistas, Knuth, Morris e Pratt, não se preocuparam em publicar o algoritmo até 1976

## Nota histórica

- Nesse meio tempo, R. S. Boyer e J. S. Moore (e, independentemente, R. W. Gosper) descobriram um algoritmo que é muito mais rápido em muitas aplicações.
- Muitos editores de texto usam esse algoritmo para busca de cadeias.
- Em 1980, M. O. Rabin e R. M. Karp desenvolveram um algoritmo tão simples quanto o de força bruta que roda virtualmente sempre em tempo proporcional a  $M + N$ .
- Além disso, o algoritmo deles estende-se facilmente a padrões bidimensionais que o torna mais útil que os outros para processamento de figuras

# Algoritmo de Força Bruta

- O método óbvio para casamento de padrão resume-se em testar, em cada posição do texto onde o padrão pode casar, se ele de fato casa
- O procedimento *Naive-String-Matcher* a seguir busca dessa maneira a primeira ocorrência do padrão P no texto T





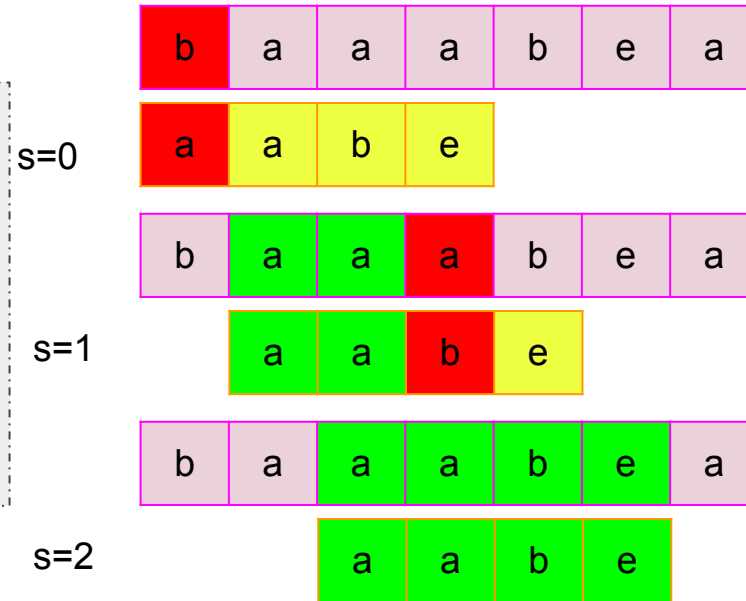
# Algoritmo Força Bruta

```

Naive-String-Matcher(P, T)
n = |T|
m = |P|
for s = 0 to n-m
    if P[1...m] == T[s+1, ...s+m]
        return s+1
return -1
  
```

Exemplo:

T	b	a	a	a	b	e	a
P	a	a	b	e			



Em  $s = 2$ , temos que  $P[1...4]$  casa com o trecho  $T[3...6]$

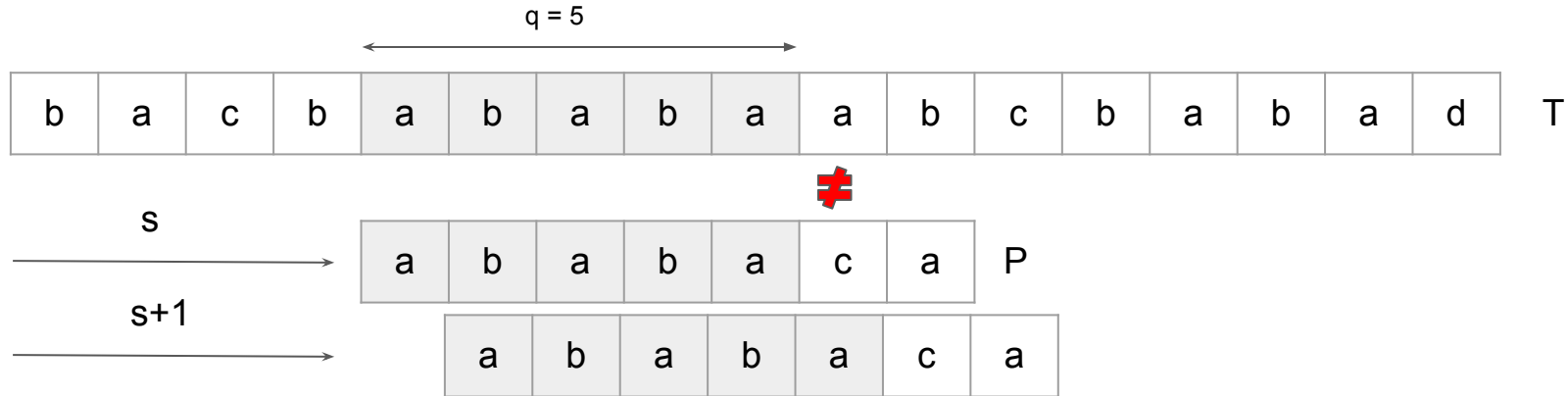
# Algoritmo de Força Bruta

- Pior caso
  - O pior caso ocorre quando, por exemplo, o padrão e o texto são os dois uma sequência de zeros seguidos por um 1
    - Por exemplo, 00001 e 00000000000000000001
  - Ou seja, quando é preciso percorrer praticamente todo P várias vezes a cada posição de T, estando P no fim da cadeia
  - Complexidade:  $m \cdot n$

# Algoritmo Knuth-Morris-Pratt (KMP)

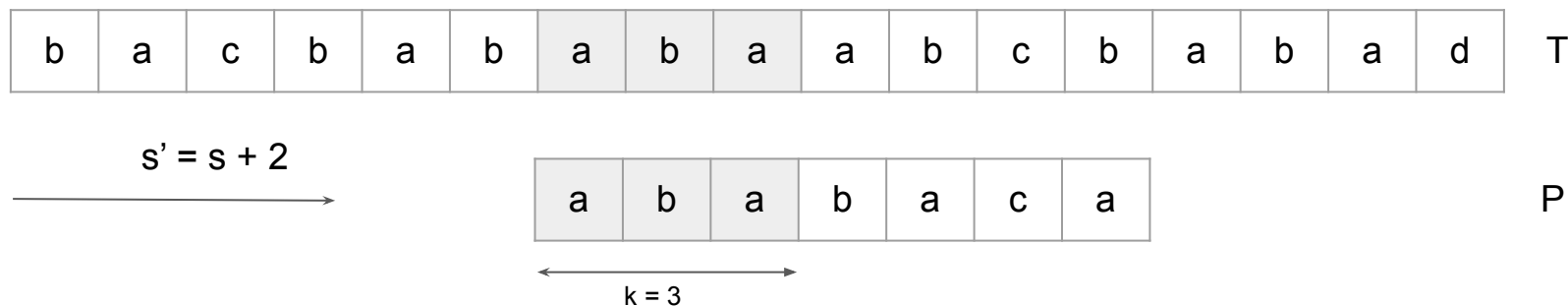
- Ideia básica:
  - No algoritmo de força bruta, quando ocorre uma diferença entre  $T[i]$  e  $P[j]$ , não seria melhor um deslocamento maior de  $P$  para a direita evitando comparações redundantes?
  - Quando ocorre um descasamento,  $P$  tem em si a informação necessária para determinar onde começar a próxima comparação, ou seja,  $P$  inclui subsequências idênticas no início de  $P$  e antes do caractere incompatível.
- É o primeiro algoritmo cujo pior caso tem complexidade de tempo linear no tamanho do texto.

# Algoritmo Knuth-Morris-Pratt (KMP) - Exemplo



Observe que o deslocamento  $s + 1$  é necessariamente não válido, uma vez que o primeiro caractere do padrão estaria alinhado com um caractere do texto que sabe-se que não casará.

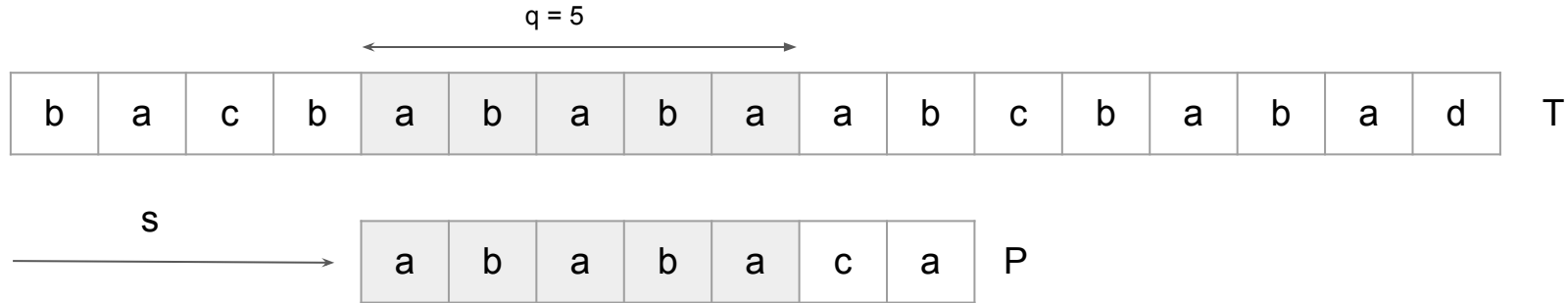
# Algoritmo Knuth-Morris-Pratt (KMP) - Exemplo



O deslocamento  $s + 2$  alinha os três primeiros caracteres do padrão  $P$  com três caracteres do texto que devem ser necessariamente correspondentes.

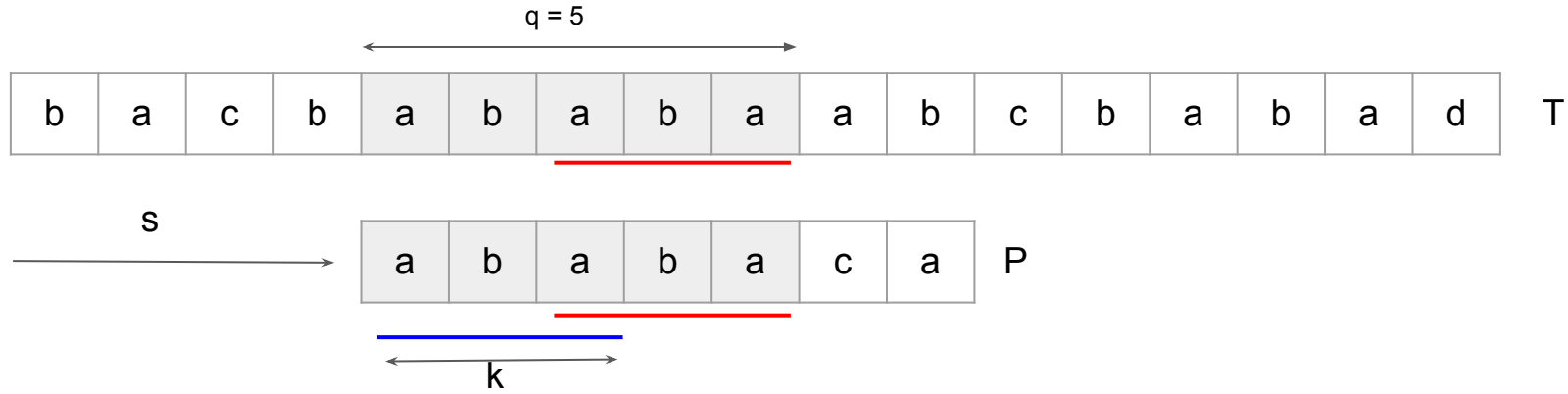
**A questão então é: temos como saber qual é o deslocamento necessário para casar esses trechos e evitar comparações desnecessárias?**

# Algoritmo Knuth-Morris-Pratt (KMP)



Sabendo que  $P[1...q]$  está contido em  $T[s+1, s+q]$ , queremos o prefixo próprio mais longo  $P[1...k]$  de  $P[1...q]$  que é também sufixo de  $T[s...q]$ . Note que, para encontrar o menor deslocamento  $s'$  equivale a encontrar o maior prefixo  $k$ , pois  $s' = s + (q - k)$ .

# Algoritmo Knuth-Morris-Pratt (KMP)



- Note que  $P[1...5]$  está contido em  $T[s+1, s+5]$  e que o **prefixo** próprio mais longo  $P[1...3]$  de  $P[1...5]$  é também **sufixo** de  $T[s+1...s+5]$ .
- Para encontrar o menor deslocamento  $s'$  equivale a encontrar o maior prefixo  $k$ , pois  $s' = s + (q - k)$ . Logo  $s' = 2$ .

# Algoritmo Knuth-Morris-Pratt (KMP)

- Tem como saber antecipadamente quantas posições devem ser deslocadas?
- Para tal, utilizamos uma função chamada de **Função Prefixo**, denotada por  $\pi$ .



# Função Prefixo

- A função prefixo encapsula o conhecimento sobre quantas posições deve-se caminhar para continuar procedendo o casamento do padrão, evitando comparações inúteis.
- Para tal, a função analisa o padrão fornecido e cria uma “tabela” de deslocamentos.
  - Para cada valor de  $q$ , armazena o número  $k$  de caracteres correspondentes do novo deslocamento  $s'$
- Formalmente:

Dado um padrão  $P[1\dots m]$ , a função prefixo para o padrão  $P$  é a função  $\pi = [1, 2, \dots, m] \rightarrow [0, 1, \dots, m-1]$  tal que  $\pi[q] = \max\{k: k < q \text{ e } p[1\dots k] \sqsubseteq P[1\dots q]\}$ .

# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$							

# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0						

# Função Prefixo

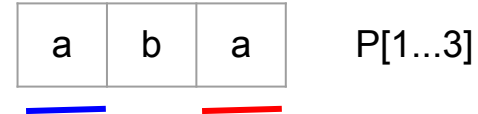
Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0					

# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

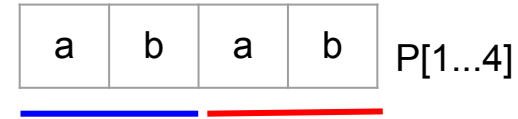
i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0	1				



# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

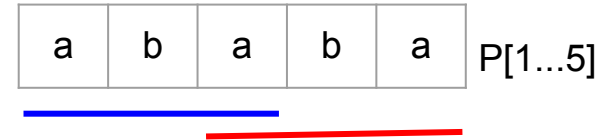
i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0	1	2			



# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0	1	2	3		



# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0	1	2	3	0	

a	b	a	b	a	c
---	---	---	---	---	---

 P[1...6]



# Função Prefixo

Exemplo do cálculo do vetor deslocamento:

i	1	2	3	4	5	6	7
P	a	b	a	b	a	c	a
$\pi$	0	0	1	2	3	0	1

a	b	a	b	a	c	a	P[1...6]
---	---	---	---	---	---	---	----------

$$\pi[q] = \max\{k: k < q \text{ e } p[1...k] \sqsupseteq P[1...q]\}$$

# Pseudocódigo KMP

1 **KMP-Matcher**( $T, P$ )

**Entrada:** Cadeias de caracteres  $T$  e  $P$

2  $n \leftarrow |T|;$

3  $m \leftarrow |P|;$

4  $\pi \leftarrow \text{ComputaFuncaoPrefixo}(P);$

5  $q \leftarrow 0;$

// considere  $i \leq n$

6 **para**  $i \leftarrow 1$  **até**  $n$  **faça**

7     **enquanto**  $q > 0$  **e**  $P[q + 1] \neq T[i]$  **faça**

8          $q \leftarrow \pi[q];$

9     **fim**

10    **se**  $P[q + 1] = T[i]$  **então**

11          $q \leftarrow q + 1;$

12    **fim**

13    **se**  $q = m$  **então**

14         **Imprima** "O padrão ocorre com deslocamento  $i - m$ ";

15          $q \leftarrow \pi[q];$

16    **fim**

17 **fim**

// Número de caracteres correspondentes

// Varre o texto da esquerda para direita

// próximo caractere não correspondente

// próximo caractere é correspondente

// P inteiro é correspondente?

// Procura próxima correspondência.

# Pseudocódigo Função Prefixo

```
1 ComputaFuncaoPrefixo( $P$ )  
   Entrada: Cadeia de caracteres  $P$   
2  $m \leftarrow |P|$ ;  
3 Crie  $\pi[1..m]$ ; é um arranjo completamente novo  
4  $\pi[1] \leftarrow 0$ ;  
5  $k \leftarrow 0$ ;  
   // considere  $q \leq m$   
6 para  $q \leftarrow 2$  até  $m$  faça  
7   enquanto  $k > 0$  e  $P[k + 1] \neq P[q]$  faça  
8      $k \leftarrow \pi[k]$ ;  
9   fim  
10  se  $P[k + 1] = P[q]$  então  
11     $k \leftarrow k + 1$ ;  
12  fim  
13   $\pi[q] \leftarrow k$ ;  
14 fim  
15 retorna ( $\pi$ );
```

## Exercício

Calcule a função prefixo para o padrão abaixo

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

# Exercício

Calcule a função prefixo para o padrão abaixo

P	a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
$\pi$																			

1 **ComputaFuncaoPrefixo**( $P$ )

**Entrada:** Cadeia de caracteres  $P$

2  $m \leftarrow |P|;$

3 **Crie**  $\pi[1..m]$ ; **é um arranjo completamente novo**

4  $\pi[1] \leftarrow 0;$

5  $k \leftarrow 0;$

Calcule a função prefixo para o padrão abaixo



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
P	a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
$\pi$	0																		

m=19

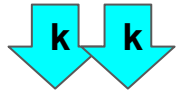
1 **ComputaFuncaoPrefixo**( $P$ )

**Entrada:** Cadeia de caracteres  $P$

- 2  $m \leftarrow |P|;$
- 3 **Crie**  $\pi[1..m]$ ; **é um arranjo completamente novo**
- 4  $\pi[1] \leftarrow 0;$
- 5  $k \leftarrow 0;$



Calcule a função prefixo para o padrão



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

P

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

m=19

$\pi$

0	0																	
---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

6 para  $q \leftarrow 2$  até  $m$  faça

7 enquanto  $k > 0$  e  $P[k+1] \neq P[q]$  faça  $k \leftarrow 0$  (F)

8 |  $k \leftarrow \pi[k];$

9 fim

10 se  $P[k+1] = P[q]$  então

11 |  $k \leftarrow k+1;$

12 fim

13  $\pi[q] \leftarrow k;$

14 fim

$P[1] == P[3]$  (V) Logo  $k \leftarrow 1$

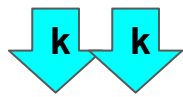
$\pi[3] \leftarrow 1$







Calcule a função prefixo para o padrão



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

P

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

m=19

$\pi$

0	0	1	2	0														
---	---	---	---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--

```

6  para q ← 2 até m faça
7      enquanto k > 0 e P[k + 1] ≠ P[q] faça
8          | k ← π[k];
9      fim
10     se P[k + 1] = P[q] então
11         | k ← k + 1;
12     fim
13     π[q] ← k;
14 fim

```

k > 0 (F)

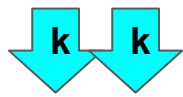
P[1] == P[6] (V) Logo k ← 1

$\pi[6] \leftarrow 1$





Calcule a função prefixo para o padrão



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

P

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

m=19

$\pi$

0	0	1	2	0	1	2	0											
---	---	---	---	---	---	---	---	--	--	--	--	--	--	--	--	--	--	--

6 para  $q \leftarrow 2$  até  $m$  faça

7 enquanto  $k > 0$  e  $P[k+1] \neq P[q]$  faça  $k > 0$  (F)

8 |  $k \leftarrow \pi[k];$

9 fim

10 se  $P[k+1] = P[q]$  então

11 |  $k \leftarrow k+1;$

12 fim

13  $\pi[q] \leftarrow k;$

14 fim

$P[1] == P[9]$  (V) logo  $k \leftarrow 1$

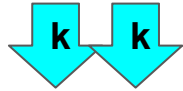
$\pi[9] \leftarrow 1$







Calcule a função prefixo para o padrão



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

P

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

m=19

$\pi$

0	0	1	2	0	1	2	0	1	2	0								
---	---	---	---	---	---	---	---	---	---	---	--	--	--	--	--	--	--	--

```

6 para q ← 2 até m faça
7   enquanto k > 0 e P[k + 1] ≠ P[q] faça
8     k ← π[k];
9   fim
10  se P[k + 1] = P[q] então
11    k ← k + 1;
12  fim
13  π[q] ← k;
14 fim

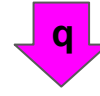
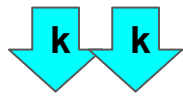
```

k > 0 (F)

P[1] == P[12] (V), logo k ← 1

π[12] ← 1

Calcule a função prefixo para o padrão



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

P

a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

m=19

$\pi$

0	0	1	2	0	1	2	0	1	2	0	1							
---	---	---	---	---	---	---	---	---	---	---	---	--	--	--	--	--	--	--

```

6 para q ← 2 até m faça
7   enquanto k > 0 e P[k + 1] ≠ P[q] faça
8     | k ← π[k];
9   fim
10  se P[k + 1] = P[q] então
11    | k ← k + 1;
12  fim
13  π[q] ← k;
14 fim

```

k == 0, k > 0(F)

P[1] == P[12] (V)!, logo k ← 1

$\pi[12] \leftarrow 1$

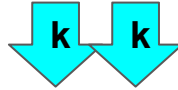








Calcule a função prefixo para o padrão



	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
P		a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
$\pi$	0	0	1	2	0	1	2	0	1	2	0	1	2	3	4	5				

m=19

```

6  para q ← 2 até m faça
7      enquanto k > 0 e P[k + 1] ≠ P[q] faça
8          |   k ← π[k];
9      fim
10     se P[k + 1] = P[q] então
11         |   k ← k + 1;
12     fim
13     π[q] ← k;
14 fim

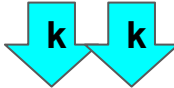


```

k > 0 e P[6] != P[17] (F)!

P[6] == P[17] (V) logo k ← 6

π[17] ← 6

Calcule a função prefixo para o padrão

																				
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
P		a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
π		0	0	1	2	0	1	2	0	1	2	0	1	2	3	4	5	6		

m=19

m=19

```

6  para q ← 2 até m faça
7      enquanto k > 0 e P[k + 1] ≠ P[q] faça
8          | k ← π[k];
9      fim
10     se P[k + 1] = P[q] então
11         | k ← k + 1;
12     fim
13     π[q] ← k;
14 fim

```

k>0 e P[7] != P[18] (F)!

P[7] == P[18] (V) logo k ← 7

π[18] ← 7





Calcule a função prefixo para o padrão

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
P		a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b	m=19
$\pi$		0	0	1	2	0	1	2	0	1	2	0	1	2	3	4	5	6	7	8	

## Exercício 2:

1. Compute a função prefixo para  $P = \text{abbarba}$ ;
2. Use o algoritmo KMP para buscar  $P$  em  $T = \text{abacaabaccabacabaabb}$ .
3. Quantas comparações o método faz?

Visualizador:

<https://cmeps-people.ok.ubc.ca/ylucet/DS/KnuthMorrisPratt.html>

# **Processamento de Cadeias de Caracteres**

Prof. José J. Camata  
camata@ice.ufjf.br